# TWEETER SUICIDAL NOTES IDENTIFICATION USING ADVERSARIAL DEEP THEORY

## BY

### MD. MAHADI RAHMAN DHRUBO
**ID: 191-15-12177**

### SADIA ISLAM
**ID: 191-15-12444**

## AND

### MAHJABIN BINTA KABIR
**ID: 191-15-12818**

This Report Presented in Partial Fulfillment of the Requirements for the Degree of Bachelor of Science in Computer Science and Engineering

Supervised By

**MD. JUEAL MIA**
Assistant Professor
Department of CSE
Daffodil International University

Co-Supervised By

**SHARUN AKTER KHUSHBU**
Lecturer
Department of CSE
Daffodil International University

# DAFFODIL INTERNATIONAL UNIVERSITY

## DHAKA, BANGLADESH
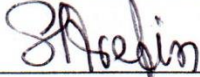
## JANUARY 2023

# APPROVAL

This Project titled **"Tweeter Suicidal Notes Identification Using Adversarial Deep Theory"**, was submitted by **Md. Mahadi Rahman Dhrubo ID No: 191-15-12177, Sadia Islam ID No: 191-15-12444, and Mahjabin Binta Kabir ID No: 191-15-12818** to the Department of Computer Science and Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of B.Sc. in Computer Science and Engineering and approved as to its style and contents. The presentation was held on 24th January 2023.

## BOARD OF EXAMINERS

**Chairman**

**Dr. Touhid Bhuiyan**
**Professor and Head**
Department of Computer Science and Engineering
Faculty of Science &amp; Information Technology
Daffodil International University

**Internal Examiner**

**Dr. Mohammad Shamsul Arefin**
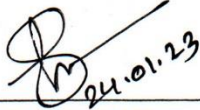**Professor**
Department of Computer Science and Engineering
Faculty of Science &amp; Information Technology
Daffodil International University

**Internal Examiner**

**Md. Sabab Zulfiker**
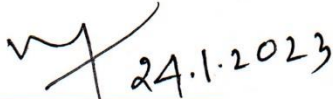**Senior Lecturer**
Department of Computer Science and Engineering
Faculty of Science &amp; Information Technology
Daffodil International University

**External Examiner**

**Dr. Ahmed Wasif Reza**
**Associate Professor**
Department of Computer Science and Engineering
East West University

# DECLARATION

We hereby declare that, this project has been done by us under the supervision of **Md. Jueal Mia, Assistant Professor, Department of CSE** Daffodil International University. We also declare that neither this project nor any part of this project has been submitted elsewhere for award of any degree or diploma.

**Supervised by:**
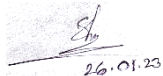
**Md. Jueal Mia**
Assistant Professor
Department of CSE
Daffodil International University

**Submitted by:**

**Md. Mahadi Rahman Dhrubo**
ID: -191-15-12177
Department of CSE
Daffodil International University

**Sadia Islam**
ID: -191-15-12444
Department of CSE
Daffodil International University

**Mahjabin Binta Kabir**
ID: -191-15-12818
Department of CSE
Daffodil International University

# ACKNOWLEDGEMENT

First we express our heartiest thanks and gratefulness to almighty God for His divine blessing makes us possible to complete the final year project/internship successfully.

We really grateful and wish our profound our indebtedness to **Md. Jueal Mia**, **Assistant Professor**, Department of CSE Daffodil International University, Dhaka. Deep Knowledge & keen interest of our supervisor in the field of "*Machine Learning*" to carry out this project. His endless patience ,scholarly guidance ,continual encouragement , constant and energetic supervision, constructive criticism , valuable advice ,reading many inferior draft and correcting them at all stage have made it possible to complete this project.

We would like to express our heartiest gratitude to **Prof. Dr. Touhid Bhuiyan, Head,** Department of CSE, for his kind help to finish our project and also to other faculty member and the staff of CSE department of Daffodil International University.

We would like to thank our entire course mate in Daffodil International University, who took part in this discuss while completing the course work.

Finally, we must acknowledge with due respect the constant support and patients of our parents.

# ABSTRACT

Depression is a common mental illness that can interfere with daily activities and productivity. Suicidal thoughts or attempts may result as a result. In today's society, suicide is a major issue. Suicide attempts should be detected and prevented early on in order to preserve people's lives. Natural Language Processing (NLP) and machine learning techniques were used to construct the platform, which was designed to interpret conversations. The proposed two-stage platform would evaluate conversation and categorize associated sentiments into four categories: "happy", "neutral", "depressive", and "suicidal". The first step of intent recognition would examine conversations and categorize associated sentiments into two categories: "YES" and "NO". We show how social media data and suicide notes can be used to identify people who are in danger of committing suicide. Suicide notes are usually written in letters and posted on websites, and they're also captured in audio and video. Suicide notes can be used as study material in NLP. This article comprehensively introduces and explores methods and algorithms from a variety of disciplines. We also investigate the implications of using the same gear and thus the security implications. For knowledge-aware suicide risk assessment, current research incorporated external knowledge utilizing knowledge bases and suicide ontology. Even this surgical technique however is actually currently there only available as well for state intervention among users whom had "checked out" either study and counseling. Still, it allows for scalable suicide risk screening, potentially identifying many people who are at risk before they engage with services in the health infrastructure. Finally, we review considering the current state of affairs work's shortcomings and offer some recommendations for further research.

# TABLE OF CONTENTS

| CONTENTS | PAGE |
|---|---|

# LIST OF TABLES

| Tables | PAGE NO |
|---|---|
| Table 1: Performance of machine learning model | 14 |
| Table 2: Result of applying algorithm accuracy score | 14 |

# LIST OF FIGURES

# CHAPTER 1

# INTRODUCTION

## 1.1 Introduction

MENTAL health disorders, such as anxiety and depression, are a growing source of concern in modern society, as they appear to be particularly severe in developed and emerging markets [1]. Suicidal detection or even suicide attempts can result from severe mental conditions that go untreated. Some internet posts contain a great deal of unpleasant material and can lead to issues like cyberstalking and cyberbullying. Consequences can be serious and dangerous, as bad information is frequently used in some type of social cruelty, resulting in gossip and even mental harm. Cyberbullying and suicide have been linked, according to research. Victims who are repeatedly exposed to unpleasant messages or situations may grow sad and desperate, and some may even commit suicide. Suicide is complicated for a variety of reasons. Suicide is more common in those who suffer from depression, however, many people who do not suffer from depression can have suicidal thoughts. Suicide variables are divided into three categories by the American Foundation for Suicide Prevention (AFSP) [1]: health, environmental, and historical factors. Suicide risk factors include mental health problems and substance abuse problems. Psychological risks such as personality and individual characteristics, cognitive variables, social factors, and negative life events were outlined by O'Connor and Nock in a comprehensive examination of the psychology of suicide.[1]Suicidal detection (SD) is a technique that uses tabular data or textual content supplied by a person to identify whether or not they have suicidal thoughts. Strange social phenomena, such as online communities agreeing on self-mutilation and copycat suicide, are natural language processing used to create an interactive platform for conducting discussions in order to recognize users' thoughts or intentions (NLP) [1][2].

The two-phased technology combines a machine learning method with Google Home mini, a virtual personal assistant, to determine suicidal intents from a person's spoken utterance (VPA) [2]. To link severe depression with suicide intentions, it's critical to find terms connected with significant negative feelings. The Google Dialogflow machine learning

technique was used to build a model to understand the user's mental or emotional states from his or her spoken speech and/or chats with the VPA [2] during the first stage intent recognition phase.

For example, [2]in 2016, a social media craze known as the "Blue Whale Game" [1] used a variety of challenges (including self-harm) to persuade game participants to commit suicide. Suicide is a major social problem that claims the lives of thousands of people each year. Almost as a response, it's crucial to acknowledge adverse outcomes and aid before sufferers conduct either themselves or. The most effective methods of preventing prospective suicide attempts are early detection and therapy. In the first stage, the model was trained using a dataset that included various sentimental phrases as well as genuine suicidal notes, which were grouped into two emotional states: 'YES' 'NO'.In the second stage, known as the emotion nurture phase, if severe depression or 'Suicidal' intent was detected, a distress signal was sent to a suicide prevention helpline; otherwise, the platform proceeded with helpful conversations. As a result, our detection method focuses on both diagnosing the problem and delivering quick resources and a treatment plan after the problem has been identified. In addition to established approaches, the use of such a novel instrument can be utilized to commence adequate resources and treatment for depression. Many difficult societal problems have been solved with AI. One of the potential applications for social good is the detection of suicide behavior using AI approaches, which should be addressed in order to meaningfully improve people's well-being. Feature selection on tabular and text data, as well as natural language representation learning, are two of the study topics. To classify suicide risks, many AI-based approaches have been used. There are still some obstacles to overcome. For training and evaluating SD, there are a restricted number of benchmarks. Statistical clues are occasionally learned by AI-powered models, although they do not always grasp people's intentions. Furthermore, many neuronal models are difficult to understand. This overview examines SD methodologies in the context of AI and machine learning, as well as specific domain applications with social implications. This article provides a thorough examination of the rapidly growing topic of SD using machine learning techniques. It gives a synopsis of current research progress as well as a forecast for future efforts. It lists certain activities as well as various data sets.

Finally, we have a debate and make some recommendations for the future. Those learning algorithms demand social media platforms involving people who were able to despair because once individuals suffering from depression, and therefore a regulated base of users might not have contemplated suicide. Through programs to inform optimization algorithms could discern between someone being at risk of self-harm or who is not, humans sought evidence of people who were just as tight as possible to someone who would suffer from mental health issues but did not even attempt themselves. Control users' social media posts develop a clear understanding against relevant background information individuals who attempt suicide can be compared. Internet users try to register with this media platform and enhance the transparency as well of their own digital counterparts, which includes social media and the internet, wearable technology, and other communications technology. Individuals must also usually fill out anyway aerial surveys which try to ask for basic demographic information etc as well as technical information about their cognitive wellness condition. They record the precise magnitude as well as periods for previous Attempts towards suicidal behavior particularly pertinent to this study. The construction analytical evaluation using predictive modeling requires social networking sites' insights by consumers before they attempted despair, in addition to a control group of most individuals who might not have attempted to assassinate yourselves. We wanted instances of people who were as near as feasible to all of those who sought would conduct themselves however could not do anything in order to teach to algorithms distinguish citizens who are already at the probability of suicidal behavior versus those that are not. Control users' social media posts develop a clear understanding of relevant resources to individuals who attempt suicide can be compared.

## 1.2 Motivation

Nowadays we are seeing frequent suicide cases in our country. Suicide seems to have become a daily practice among everyone. Some are committing suicide due to their family problems, some are committing suicide due to not getting a job, some are committing suicide due to failure in love. These suicides are a threat to our country today, these suicide cases have become a cause of concern for all of us. Nowadays, before suicide, people post

on their social accounts why they are going to commit suicide or sometimes post in such a way that it is not clear that they are actually going to commit suicide. It is mainly because of all these factors that we are motivated to research this topic and we are continuing this research work.

## 1.3 Objective

We show how social media data and suicide notes can be used to identify people who are in danger of committing suicide. Suicide notes are usually written in letters and posted on websites, and they're also captured in audio and video. Suicide notes can be used as study material in NLP. This article comprehensively introduces and explores methods and algorithms from a variety of disciplines. We also investigate the implications of using the same gear and thus the security implications.

## 1.4 Expected Outcome

After determining whether or not a person is acting suicidally, sentiment analysis must be properly analyzed and evaluated. The anticipated outcome would be more error-free after training and retraining the data with the proper method. We can contrast real positive and false positive alarms to better understand how this automation system is being evaluated [18]. We can anticipate that 4% to 8% (let's estimate 6%) of a population of 1000 people being screened by this intelligence system would have a plan to commit suicide. This suggests that out of 1000 persons, 60 would try suicide while the other 940 would not. If we assume that our system has a 10% false alert rate, then 144 persons will be classified as "suicidal," of which about 35% will really kill themselves. The false-positive label would decrease as the false alert frequency decreased. If the system flagged the suspect false positive warnings, it would be more useful.

# CHAPTER 2

# LITERATURE REVIEW

## 2.1 Introduction

By using tabular information on a person or text that was written by a person, suicidal ideation detection (SID) ascertains whether the person has suicidal ideas or thoughts. An rising number of people are choosing to engage with others online due to advancements in social media and online anonymity. Online communication is evolving into a new medium for people to communicate their emotions, pain, and suicidal thoughts. As a result, social media content mining can help with suicide prevention. Online channels have thus started to operate as a surveillance tool for suicidal thinking. Potential suicide victims may act out their thoughts of killing themselves in role-playing, fleeting thoughts, and suicide plots. Before a catastrophe occurs, SID should identify these dangers in intents or behaviors.

## 2.2 Paper Review

A number of domains, including computer vision, natural language processing, and medical diagnosis, have found great success using deep learning. Researchers in the area of suicide prevention and automatic SD use it extensively. Feature engineering is not required to understand the features of a given piece of text. Some individuals are using deep neural networks to process the retrieved characteristics at the same time as others. Research on SD utilizing multimodal data has benefited from machine learning's growing popularity, and this technology offers a promising alternative for effective early warning systems. Text-based approaches for feature extraction and deep learning are the focus of current research. For example, deep learning models like CNN[1][21] and LSTM [20][21] are often utilized by academics in conjunction with conventional NLP characteristics like the topical, syntactic, emotional, and readability aspects of the TF-IDF [1]. These techniques, notably DNNs[1] with autonomous feature learning, enhanced prediction performance and early success in understanding suicidal intention.

On the contrary, certain techniques may just provide statistical evidence and a lack of common sense. Knowledge bases and a suicide ontology were included into current

research for knowledge-aware suicide risk assessment. Knowledge-based detection has taken a huge step forward. Online user content, electronic health records, suicide diaries, and other forms of survey data collection comprise the four primary SD application types. The use of questionnaires and electronic health records (EHRs) relies significantly on social workers and mental health experts, and requires self-report assessment or patient-clinician encounters to be made. Suicide notes are written by those who intend to kill themselves before they do so. They're commonly expressed in the form of letters, but they may also be recorded in audio and video form and shared online. Using suicide notes as study material in NLP is a good idea [21]. There has been a lot of study done on suicide detection in Chinese microblogs. Content analysis and study on suicide risk variables will still benefit from their use. The final online user content domain may be one of the most promising early warning and suicide prevention tools when paired with machine learning technologies. As digital technology advances at a fast pace, user-generated content will become more important in SD. When it comes to monitoring suicide risk in the future, wearable gadgets' health data may be quite helpful.

## 2.3 Summary of Research

Four areas make up the majority of SID applications: questionnaires, electronic health records, suicide notes, and online user material.

A summary of the categories, data sources, and methodology is provided in Table II. Among these four key areas, questionnaires and electronic health records (EHRs) rely heavily on social workers or other mental health professionals and demand self-report measurement or patient-clinician contacts. Because many people who attempt suicide do so shortly after writing suicide notes, suicide notes have limited potential for early prevention. They do, however, offer a useful resource for content analysis and the investigation of suicide-related issues.

When combined with machine learning approaches, the final online user content domain offers one of the most promising means of early detection and suicide prevention.

# CHAPTER 3

# METHODOLOGY

Kaggle provided us with the data set that we used in our article. We have contributed over 100 pieces of data to this data collection in addition to the Kaggle data. The purpose of our data set is to detect suicidal tendencies. Determine whether a social media post or letter is suicidal by looking at the content. This data set is divided into two categories. "Yes" and "No".



Figure 1: Full Process of our model building. We collect data then process our data and then build the model.

## 3.1 Training of ML/DL/NLP Algorithms

Written text provided was supplied since parameters into Google Colaboratory ML algorithms [5] and python programming deep learning methods utilizing kaggle to take advantage of better computational capacity. Humans made do with a network computing environment with the following specs in Google Colaboratory [5]. These ML algorithms were trained using the Sklearn library. As a comparison to NLP algorithms, we performed

tests with several ML algorithm implementations provided by Scikit-learn. With the grid search technique, they practiced several

as well as classifications and fine-tuned their scalability and performance can optimize these algorithms' efficiency. A technique like this thoroughly tests several combinations of hyperparameters in order to discover the optimum one. In the end, the models that made use of the stochastic gradient descent, logistic regression, support vector machine, and random forest classifier methods generated the best overall performance results. Number of occurrences of No and Yes: (a) without data balancing (the original dataset); (b) with 80 percent used for training and validation; and (c) with 20 percent employed for testing.

## 3.2 Data Preparation

The texts were subjected to the following methods in order to prepare the data: 1. Standardized asset finders [21], Characters, emails, & numerals are examples of terms that should be removed from the text.[5] 2. Stop words[5]: "as," "e," "os," "de," "para," "com," "sem," are all examples of "stop words," which are words that don't aid in the analysis. 3. Tokenization [5][21] is a process for breaking down texts into smaller pieces called tokens (the phrase is broken down into words). 4. Stemming is a technique for reducing accent on our tongue most incidental component[5]. It's a good idea to use the same word for all of the following phrases: "cat," "gata," "gatos," and "gata," which will all be shortened to "cat" (the stem). 5. Frequency–Inverse Document Frequency[5][21] Word scoring can be greatly improved by using (TF-IDF): a statistical technique that evaluates how essential a term is to a phrase in a collection of sentences.

## 3.3. Supervised & Unsupervised Learning Process

**3.3.1. Stochastic Gradient Descent(SGD):** The stochastic nature of a system or process refers to its potential to have an unexpected outcome. Random samples from the whole data set are chosen for each iteration of Stochastic Gradient Descent. Gradient Descent refers to the total number of samples utilized in each iteration to calculate the gradient as a "batch." In normal Gradient Descent optimizations, such as Batch Gradient Descent, the batch is considered to be the complete dataset. The problem arises when our datasets are big enough to benefit from utilizing the complete dataset in order to reach the minima in a less noisy and random way. Traditional Gradient Descent optimization requires you to

utilize every sample in your dataset to complete one iteration of the algorithm, and you'll have to do this for every iteration until the minima are achieved. Thus, computation becomes prohibitively costly. This problem may be solved using Stochastic Gradient Descent. When using SGD, each iteration is carried out with a single sample and a single batch size. The sample is mixed and chosen at random to execute the iteration.

**3.3.2. Logistic Regression:** Logistic regression is one of the most widely used Machine Learning algorithms, and it is used within the context of the Supervised Learning methodology. It is a technique for determining the likelihood of a category outcome based on the values of a number of independent factors. Logistic regression is used to predict the result of an analysis using a categorical dependent variable. Because of this, the value that is produced must be either discrete or categorical. It is possible for it to be Yes or No, 0 or 1, true or untrue, and so on; but, rather of providing precise values such as 0 and 1, it provides probabilistic readings are somewhere between 0 and 1. Logistic Regression is quite similar to Linear Regression, with the exception being the way in which they are used. When dealing with regression issues, the linear regression method is used, while the logistic regression approach is utilized when dealing with classification issues. An S-shaped logistic function is used instead of a regression line to accurately forecast the two highest possible values in a logistic regression equation (0 or 1).

Using the logistic function, we can predict whether or not a cell is malignant, whether or not a mouse is fat, and many other things. Logistic regression is an important machine learning approach because it can utilize both discrete and continuous datasets to provide probabilities and classify incoming data.

**3.3.3. SVM (Support Vector Machine):** It is a kind of machine learning known as a Support Vector Machine (SVM). Although [21] is typically viewed of as a classification approach, it may also be used to handle regression issues. It's capable of handling both continuous and categorical data with ease. SVM [20] builds a hyperplane in multidimensional space in order to distinguish between distinct classes. The optimal hyperplane is generated repeatedly using SVM and then used to minimize an error. As stated before, SVM's primary purpose is to identify a maximum marginal hyperplane (MMH) that divides a dataset into classes as equally as feasible. The SVM algorithm is a

fantastic classifier. It's a supervised learning method that's mostly used to classify data into several groups and subcategories. It is necessary to train SVM using label data. Classification and regression problems may both be addressed with SVM. SVM uses a decision boundary, a hyperplane between two classes, to split or categorize them. SVM may also be used to classify images and identify objects.

**3.3.4. Random Forest:** As a machine learning algorithm, random forest [22] may be used for classification and regression tasks. For classification, it uses a majority vote, whereas for regression, it uses an average. The Random Forest Algorithm's ability to work with data sets including both continuous and categorical variables, as in regression and classification, is a critical feature. Compared to other categorization methods, this one is better. In order to understand how the random forest works, we must first learn about ensemble techniques. To put it simply, a "ensemble" is a collection of individuals working together to combine different methodologies. A better approach is to make use of a collection of patterns rather than only depending on the notion when forecasting. The ensemble employs two distinct methods:

**3.3.4.1. Bagging:** That produces something different chunk of learning with replacement its outcome was selected via popular election on representative data sets voting. Consider the random decision forests algorithm.

**3.3.4.2. Boosting:** It transforms via transforming special needs students into excellent scholars building sequential models with the utmost precision. AdaBoost and XGBOOST are two examples.

**3.4. Feature Engineering**

We transfer feature engineering to text data based on our comprehension. First and foremost, it is critical that we fully comprehend our difficulties. Some techniques are only applicable to a specific sort of problem. For example, we may need to extract grammatical aspects from data for some situations, but for others, we may merely need to retrieve the most frequently occurring words. It's vital to remember that feature engineering differs from the other forms of data in NLP. We're working with language or texts in NLP, thus

to get inputs for our machine learning models, we'd have to convert our text into some kind of numeric representation that computers can understand. One of our objectives is to render our text in a computer-friendly format.

**3.4.1. Parsing:** Parsing is the practice of breaking down a sentence (or other content) into smaller bits to better grasp its syntactic structure and meaning. To analyze sentences in NLP, rules of context-free grammar (CFG) or probabilistic context-free grammar (PCFG) are utilized. Building a parser from the ground up is a difficult endeavor in and of itself. We'd choose a grammar, such as CFG or PCFG, and then decide what kind of parser to make. We'd then use that information to create our parser by implementing certain algorithms. We don't have to do it every time and can make use of pre-made tools.

### 3.5. Statistical Methods

This is a more advanced feature extraction method that extracts characteristics from text data using statistics and probability principles. Now we'll take a look at the most often used statistical feature extraction method.

**3.5.1. Term Frequency-Inverse Document Frequency (TF-IDF):** Words that appear frequently in a paper are more important than those that appear only once or twice. However, words like "a," "the," and others that appear repeatedly are unimportant and have no significance. The TF-IDF [1][20]-[22] seeks to meet these two requirements by assisting us in extracting meaningful words from documents. The term frequency method is the first component of TF-IDF. Term-Frequency estimates the frequency of each term in a document/dataset, as the name suggests. The following is the formula for determining the TF of a term t.

Let us now discuss the IDF component. IDF gives terms that only appear in a few documents with more weight. Such phrases are helpful in distinguishing between various texts. We normally don't have to write the TF-IDF algorithm ourselves and may instead utilize sci-kit-learn.

## 3.6. Advanced Methods

Because they try to map a word, sentence, or text to a fixed-length vector of real numbers, these approaches are sometimes known as vectorized methods. The purpose of this method is to derive lexical and distributional semantics from a piece of text. The meaning expressed by the words is referred to as lexical semantics, but distributional semantics refers to discovering meaning based on various distributions in a corpus.

### 3.6.1. Word2Vec

Word2vec is a set of connected models that may be used to generate what are known as word embeddings [20], [22]. These are two-layer neural networks that are quite shallow and have been trained to reconstruct the linguistic contexts of words. After training, word2vec models are able to map each word to a vector consisting of several hundred components that describe that word's link to other words. This vector is intended to serve as a representation of the neural network's hidden layer. Word2vec[22] may construct neural word embeddings by using either skip-grams or a continuous bag of words as its input (CBOW). A research team at Google lead by Tomas Mikolov was responsible for developing it. After that time, many researchers have researched and provided explanations of the method.

### 3.6.2. Word Cloud

Using a simple but powerful visual representation object, a word cloud [5][20] shows the most often used terms in a document in bigger, bolder letters and in a variety of colors. The smaller a phrase is, the less important it is.

Figure 2: Word Cloud Visualization of first 2000 text data. And here text length is 2202856.

### 3.7. Evaluation Metrics & Results

The following recommendations have been used to evaluate the various forms of segmentation and classification techniques:

1) Precision = TruePositives(tp)/TruePositives(tp) + FalsePositives(fp) ……………………….[20]
2) Recall = TruePositives(tp) / TruePositives(tp) + FalseNegatives(fn) ………………………[20]
3) F1 Score = 2*((precision*recall)/(precision+recall)) ….[20]
4) Accuracy = tp + tn / tp + tn + fp + fn ……..[20]

True positives are counted as tp, false positives as fp, and false negatives as tn, whereas false positives are counted as fp and false negatives as fn.

If we measure the performance of the four algorithms that we have applied in our paper, we can see that logistic regression and XGBoost's recall give better results than svm and random forest. If we talk about precision then logistic regression and random forest give good results. In the four methods we've tested, the logistic regression and random forest

models both performed well. Regression models based on logistic data and random forest data have both shown good results so far. Below it is highlighted through table 1:

Table 1.  Performance of machine learning model

| Algorithm | Accuracy | Precision | Recall | F-Score |
|---|---|---|---|---|
| Logistic Regression | 0.72 | 0.73 | 0.78 | 0.70 |
| XGBoost | 0.70 | 0.71 | 0.76 | 0.68 |
| SVM | 0.70 | 0.71 | 0.73 | 0.69 |
| Random Forest | 0.71 | 0.72 | 0.74 | 0.70 |

Here, below table 2 showing algorithms accuracy score:

Table 2.  Result of applying algorithm accuracy score

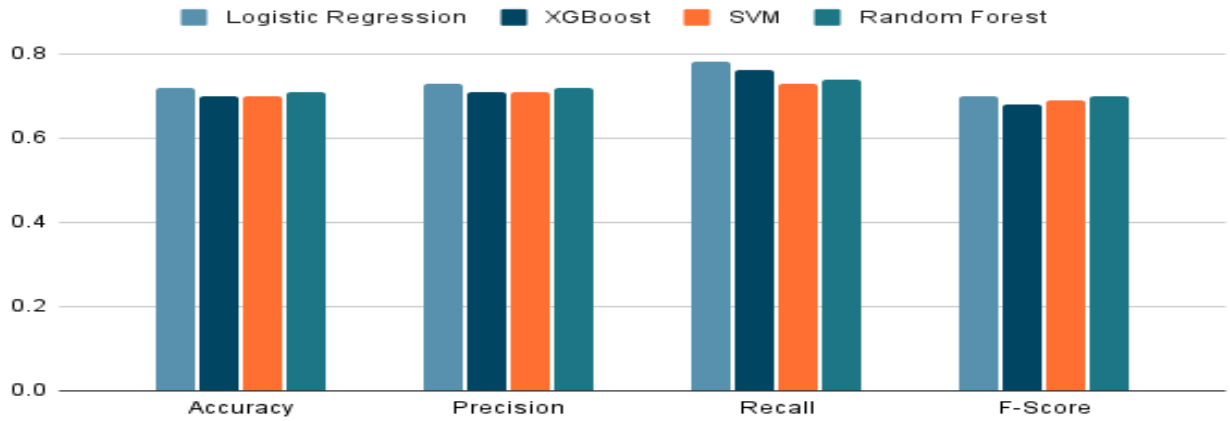| Algorithm | Accuracy |
|---|---|
| Logistic Regression | 0.72 |
| XGBoost | 0.70 |
| SVM | 0.70 |
| Random Forest | 0.71 |

Figure 5: Logistic Regression, XGBoost, SVM, Random Forest Algorithm Accuracy, Precision, Recall, F-Score.

Figure 5 above shows the precision, recall, f1 score, and accuracy plot of each algorithm using histogram.

# CHAPTER 4

# CONCLUSION

Suicide prevention is still a critical task in today's culture. Early diagnosis of suicidal detection is a crucial and effective method of suicide prevention. A novel multi-platform methodology for identifying user intents in depressed people was deployed in this study. Speech content analysis was used as an objective tool for diagnosing depression and suicidal thoughts. As a result, the methodology proposed in this study produced application-based findings and worked pretty well in recognizing users' varied emotional intentions, with an overall accuracy of 70% and accuracies of 68-72 percent in detecting Depressive and Suicidal intents. The machine learning techniques were effective far enough to be immensely useful through an originally envisioned surveillance system, but the platform's constituents were not operationally feasible. Even though there are various reasons for thinking about harnessing evidence anything from these trading algorithms to help facilitate external intervention, the majority of the general public has strongly indicated organized opposition to incomparable concerted efforts. Future plans, online social material is extremely likely to be the primary medium for SD. As a result, new methods for detecting online texts containing suicidal intent must be developed to bridge the gap between professional mental health identification and automatic machine detection in the hopes of preventing suicide. The diagnostic test can also be recommended for someone who is depressed and used at home to provide informal, ongoing support. As a result, the proposed model can be used in a variety of ways, whether alone or in conjunction with mental health services.

# REFERENCES

[1] Ji, S., Pan, S., Li, X., Cambria, E., Long, G. and Huang, Z., 2020. Suicidal ideation detection: A review of machine learning methods and applications. IEEE Transactions on Computational Social Systems, 8(1), pp.214-226.

[2] Hassan, S.B., Hassan, S.B. and Zakia, U., 2020, November. Recognizing Suicidal Intent in Depressed Population using NLP: A Pilot Study. In 2020 11th IEEE Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON) (pp. 0121-0128). IEEE.

[3] Coppersmith, G., Leary, R., Crutchley, P., and Fine, A., 2018. Natural language processing of social media as screening for suicide risk. Biomedical informatics insights, 10, p.1178222618792860.

[4] Sahu, S., Ramachandran, A., Gadwe, A., Poddar, D. and Satavalekar, S., 2021. Detection of Depression and Suicidal Tendency Using Twitter Posts. In Emerging Technologies in Data Mining and Information Security (pp. 767-775). Springer, Singapore.

[5] Diniz, E.J., Fontenele, J.E., de Oliveira, A.C., Bastos, V.H., Teixeira, S., Rabêlo, R.L., Calçada, D.B., Dos Santos, R.M., de Oliveira, A.K. and Teles, A.S., 2022, April. Boamente: A Natural Language Processing-Based Digital Phenotyping Tool for Smart Monitoring of Suicidal Ideation. In Healthcare (Vol. 10, No. 4, p. 698). MDPI.

[6] Tadesse, M.M., Lin, H., Xu, B. and Yang, L., 2019. Detection of suicide ideation in social media forums using deep learning. Algorithms, 13(1), p.7.

[7] Guidère, M., 2020, March. NLP Applied to Online Suicide Intention Detection. In HealTAC 2020.

[8] Tadesse, M.M., Lin, H., Xu, B. and Yang, L., 2019. Detection of depression-related posts in reddit social media forum. IEEE Access, 7, pp.44883-44893.

[9] Mulholland, M. and Quinn, J., 2013, October. Suicidal tendencies: The automatic classification of suicidal and non-suicidal lyricists using nlp. In Proceedings of the sixth international joint conference on natural language processing (pp. 680-684).

[10] Litvinova, T.A., Seredin, P.V., Litvinova, O.A. and Romanchenko, O.V., 2017. Identification of suicidal tendencies of individuals based on the quantitative analysis of their internet texts. Computación y Sistemas, 21(2), pp.243-252.

[11] Reddy, P.P., Suresh, C., Rao, V.K., Chandana, K.S., Sowkya, S. and Akhila, R., 2021. Vocal Analysis to Predict Suicide Tendency. In Proceedings of International Conference on Advances in Computer Engineering and Communication Systems (pp. 481-488). Springer, Singapore.

[12] Yatapala, K.Y.D.H.T. and Kumara, B.T.G.S., 2021, December. Detection of Suicide Ideation in Twitter using ANN. In 2021 6th International Conference on Information Technology Research (ICITR) (pp. 1-5). IEEE.

[13] Ghosh, M., Chatterjee, S., Das, D., Chanda, K. and Roy, S., An Analytical Study for Predicting Suicidal Tendencies Using Machine Learning Algorithms.

[14] Chandra, S., Bhattacharya, S. and Kundu, S., 2021. Suicide Ideation Detection in Online Social Networks: A Comparative Review. In Proceedings of International Conference on Innovations in Software Architecture and Computational Systems (pp. 151-167). Springer, Singapore.

[15] Chanda, K., Ghosh, A., Dey, S., Bose, R., and Roy, S., 2022. Smart Self-Immolation Prediction Techniques: An Analytical Study for Predicting Suicidal Tendencies Using Machine Learning Algorithms. In Smart IoT for Research and Industry (pp. 69-91). Springer, Cham.

[16] Dasari, N., Mittapalli, R., Dhule, A. and Wagh, S., PREDICTION OF SUICIDE USING SOCIAL MEDIA DATA.

[17] Ao, Z., Lai, J., Lu, W. and Mo, H., 2021, September. Comparisons of LSTM and other Machine Learning Approaches on Predicting Suicidal Tendency Regarding to Social

Media Posts. In 2021 International Conference on Computers and Automation (CompAuto) (pp. 19-23). IEEE.

[18] Dhelim, S., Chen, L., Ning, H., and Nugent, C., 2022. Artificial Intelligence for Suicide Assessment using Audiovisual Cues: A Review. arXiv preprint arXiv:2201.09130.

[19] Schoene, A.M., Turner, A., De Mel, G.R. and Dethlefs, N., 2021. Hierarchical Multiscale Recurrent Neural Networks for Detecting Suicide Notes. IEEE Transactions on Affective Computing.

[20] Valeriano, K., Condori-Larico, A. and Sulla-Torres, J., 2020. Detection of suicidal intent in Spanish language social networks using machine learning. International Journal of Advanced Computer Science and Applications, 11, pp.688-695.

[21] Tadesse, M.M., Lin, H., Xu, B. and Yang, L., 2019. Detection of suicide ideation in social media forums using deep learning. Algorithms, 13(1), p.7.

[22] NARYNOV, S., MUKHTARKHANULY, D., KERIMOV, I. and OMAROV, B., 2019. Comparative analysis of supervised and unsupervised learning algorithms for online user content suicidal ideation detection. Journal of Theoretical and Applied Information Technology, 97(22), pp.3304-3317.

[23] M. Guidère, 'NLP Applied to Online Suicide Intention Detection,' HealTAC 2020, Mar. 2020. [24] Q. Zhong, E.W. Karlson, B. Gelaye, and et al., 'Screening pregnant women for suicidal behavior in electronic medical records: diagnostic codes vs. clinical notes processed by natural language processing,' BMC Med. Info. Decis. Mak., vol.18, no. 30, May 2018. [25] K. Haerian, H. Salmasian and C. Friedman, 'Methods for identifying suicide or suicidal ideation in EHRs,' Annual Symposium proceedings, AMIA Symposium, vol. 2012, pp. 1244–1253, Nov. 2012.

[26] Intent matching - google clouds. Online. Available https://cloud.google.com/dialogflow/es/docs/intents-matching.

[27] Twilio Conversations API. Online. Available: https://www.twilio.com/conversations.

[28] Traning Dataset. [Online]. Available: https://www.kaggle.com/ritresearch/happydb.

[29] Traning Dataset. [Online]. Available: https://github.com/CodeWritingCow/suicide-notes. [30] M. Canonico, L. De Russis, 'A Comparison and Critique of Natural Language Understanding Tools,' Ninth International Conference on Cloud Computing, GRIDs, and Virtualization. ((Intervento presentato al convegno CLOUD COMPUTING, pp. 110-115, Barcelona, Spain, Feb. 18–22, 2018.

[31] M. McShane, 'Natural Language Understanding (NLU, not NLP) in Cognitive Systems,' AI Magazine, vol. 38, no. 4, pp. 43–56, 2017. https://doi.org/10.1609/aimag.v38i4.2745.

[32] Low Code Dialogflow bots. Mining Business Data. [Online]. Available: https://miningbusinessdata.com.

[33] A. Bradley, and P. Andrew, 'The use of the ares under the ROC curve in the evaluation of machine learning algorithms,' Pattern Recognition, vol. 30, no. 7, pp. 1145–1159, 1997.

[34] R. Silipo, and Maari. Widmann, 'Confusion Matrix and Class Statistics, towards data science,'. [Online]. Available: https://towardsdatascience.com/confusion-matrix-and-class-statistics68b79f4f510b.

[35] G. Coppersmith, K. Ngo, R. Leary, and A. Wood, "Exploratory analysis of social media prior to a suicide attempt," in Proc. 3rd Workshop Comput. Lingusitics Clin. Psychol., 2016, pp. 106–117.

[36] P. Solano et al., "A Google-based approach for monitoring suicide risk," Psychiatry Res., vol. 246, pp. 581–586, Dec. 2016.

[37] H. Y. Huang and M. Bashir, "Online community and suicide prevention: Investigating the linguistic cues and reply bias," in Proc. CHI Conf. Hum. Factors Comput. Syst., 2016, pp. 1–5.

[38] M. De Choudhury and E. Kıcıman, "The language of social support in social media and its effect on suicidal ideation risk," in Proc. 11th Int. AAAI Conf. Web Social Media, 2017, p. 32.

 [39] M. E. Larsen et al., "The use of technology in suicide prevention," in Proc. 37th Annu. Int. Conf. IEEE Eng. Med. Biol. Soc. (EMBC), Aug. 2015, pp. 7316–7319.

[40] N. Masuda, I. Kurahashi, and H. Onari, "Suicide ideation of individuals in online social networks," PLoS ONE, vol. 8, no. 4, Apr. 2013, Art. no. e62262.

[41] S. Chattopadhyay, "A study on suicidal risk analysis," in Proc. 9th Int. Conf. e-Health Netw., Appl. Services, Jun. 2007, pp. 74–78.

[42] D. Delgado-Gomez, H. Blasco-Fontecilla, F. Sukno, M. Socorro Ramos-Plasencia, and E. Baca-Garcia, "Suicide attempters classification: Toward predictive models of suicidal behavior," Neurocomputing, vol. 92, pp. 3–8, Sep. 2012.

[43] S. Chattopadhyay, "A mathematical model of suicidal-intent-estimation in adults," Amer. J. Biomed. Eng., vol. 2, no. 6, pp. 251–262, Jan. 2013.

 [44] W. Wang, L. Chen, M. Tan, S. Wang, and A. P. Sheth, "Discovering fine-grained sentiment in suicide notes," Biomed. Informat. Insights, vol. 5, no. 1, p. 137, 2012.

[45] A. Abboute, Y. Boudjeriou, G. Entringer, J. Azé, S. Bringay, and P. Poncelet, "Mining Twitter for suicide prevention," in Proc. Int. Conf. Appl. Natural Lang. Data Bases/Inf. Syst. Cham, Switzerland: Springer, 2014, pp. 250–253.

[46] E. Okhapkina, V. Okhapkin, and O. Kazarin, "Adaptation of information retrieval methods for identifying of destructive informational influence in social networks," in Proc. 31st Int. Conf. Adv. Inf. Netw. Appl. Workshops (WAINA), Mar. 2017, pp. 87–92.

[47] M. Mulholland and J. Quinn, "Suicidal tendencies: The automatic classification of suicidal and non-suicidal lyricists using NLP," in Proc. IJCNLP, 2013, pp. 680–684.

[48] X. Huang, L. Zhang, D. Chiu, T. Liu, X. Li, and T. Zhu, "Detecting suicidal ideation in chinese microblogs with psychological lexicons," in Proc. IEEE 11th Int. Conf.

Ubiquitous Intell. Comput. Autonomic Trusted Comput. IEEE 14th Int. Conf. Scalable Comput. Commun. Associated Workshops, Dec. 2014, pp. 844–849.

[49] X. Huang, X. Li, T. Liu, D. Chiu, T. Zhu, and L. Zhang, "Topic model for identifying suicidal ideation in chinese microblog," in Proc. 29th Pacific Asia Conf. Lang., Inf. Comput., 2015, pp. 553–562.

# PLAGIARISM